

HW#5

Jiale Shi

1 Factor analysis

A. One problem with mixture models is that they only use a single latent variable to generate the observations. One model that addresses this issue is the Factor analysis. FA is a low rank parameterization of an MVN. In this context,

- (a) Derive an expression for the number of independent parameters in the factor analysis model.
 - (b) Show that the factor analysis model is invariant under rotations of the latent space coordinates.
- B. We consider a data set of $D = 11$ variables and $N = 387$ cases describing various aspects of cars, such as the engine size, the number of cylinders, the miles per gallon (MPG), the price, etc. Fit a $L = 2$ dimensional factor analysis model. Plot the scores in \mathbb{R}^2 to visualize the results. To get a better understanding of the “meaning” of the latent factors, project unit vectors corresponding to each of the feature dimensions $e_1 = (1, 0, \dots, 0)$, $e_2 = (0, 1, \dots, 0)$ etc into the low dimensional space. Represent using **biplot**. Provide your discussion.

The required data set can be downloaded from [this link](#).

A. FA can be thought of as away of specifying a joint density model on X using a small number of parameters.

$$\begin{aligned} P(x_i | \theta) &= \int N(x_i | Wz_i + \mu, \gamma) N(z_i | \mu_0, \Sigma_0) dz_i \\ &= N(x_i | W\mu_0 + \mu, \gamma + W\Sigma_0 W^T) \end{aligned}$$

We can set $\mu_0 = 0$ without a loss of generality

We can absorb $W\mu_0$ into μ , $W\mu_0 + \mu \rightarrow \mu$

We can set $\Sigma_0 = I$

We can emulate a correlated prior by defining a new weight matrix $\tilde{W} = W\Sigma^{\frac{1}{2}}$

$$\begin{aligned}\text{cov}[x|\theta] &= W \Sigma_0 W^T + \Psi \\ &= (\tilde{W} \Sigma_0^{-\frac{1}{2}}) \Sigma_0 (\tilde{W} \Sigma_0^{-\frac{1}{2}})^T + \Psi \\ &= \tilde{W} \tilde{W}^T + \Psi\end{aligned}$$

We thus see that FA approximates the covariance matrix of visible vector using a low-rank decomposition

$$C \triangleq \text{cov}[x] = WW^T + \Psi$$

W only uses $O(MD)$ parameters, which allows a flexible compromise between a full covariance Gaussian with $O(D^2)$ and a diagonal covariance with $O(D)$ parameters.

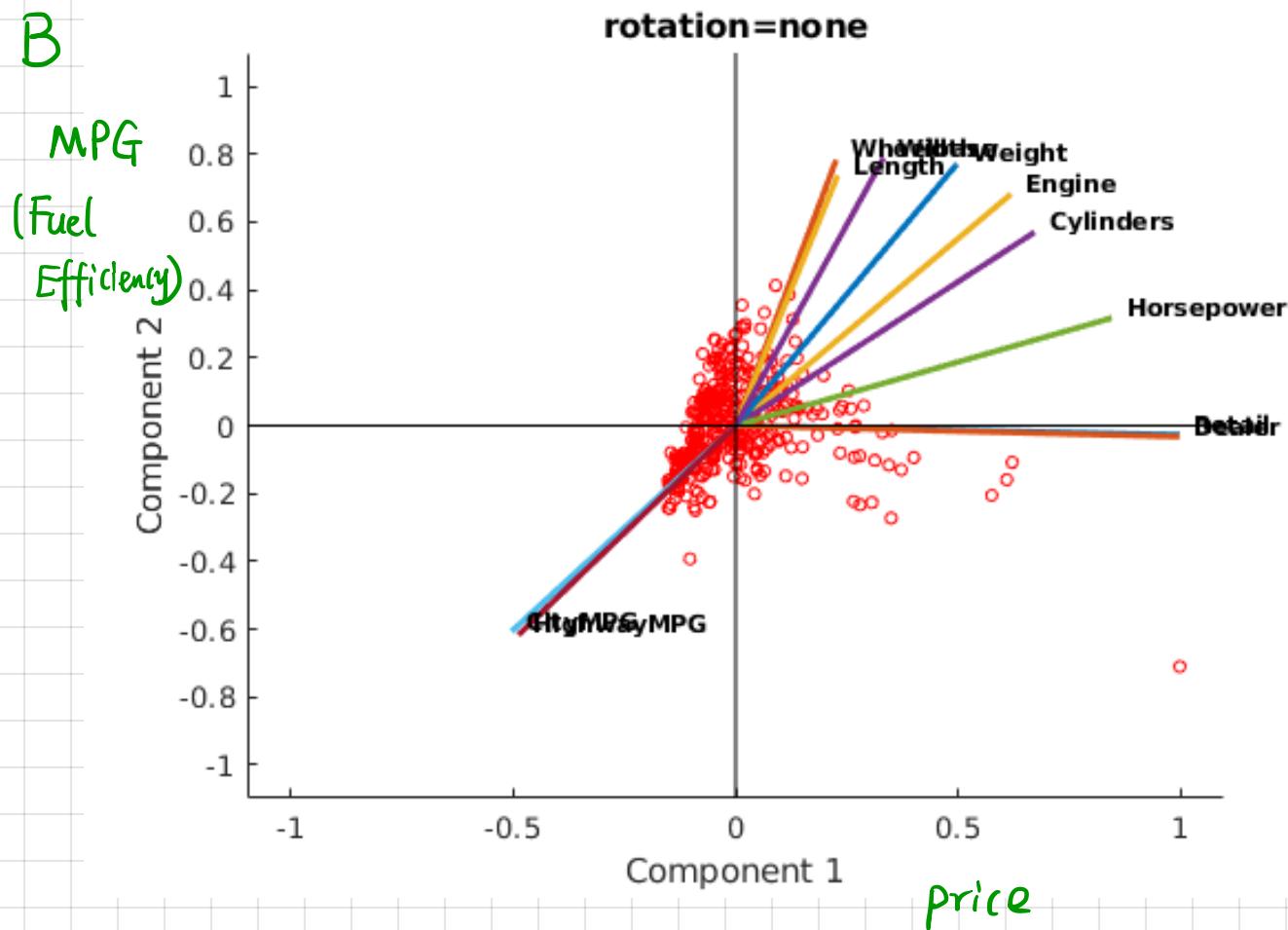
(b).

Consider an arbitrary orthogonal rotation matrix R , satisfying $RR^T = I$. Let us define $\tilde{W} = WR$

$$\tilde{W} \tilde{W}^T = WR R^T W^T = WW^T$$

We can immediately see that the likelihood $p(x_i|\theta) = \mathcal{N}(x_i|\mu, \Psi + WW^T)$ remains the same.

B



Discussion:

the horizontal axis represents price, corresponding to the features labeled "dealer" and "retail", with expensive cars on the right.

The vertical axis represents MPG, fuel efficiency versus size: heavy vehicles are less efficient and are higher up, whereas light vehicles are more efficient and are lower down.

2 PCA and KPCA

A. Consider a linear-Gaussian latent-variable model having a latent space distribution $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ and a conditional distribution for the observed variable $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Phi})$ where $\boldsymbol{\Phi}$ is an arbitrary symmetric, positive definite noise covariance matrix. Now suppose that we make a non-singular linear transformation of the data variables $\mathbf{x} \rightarrow \mathbf{Ax}$, where \mathbf{A} is a $D \times D$ matrix. If $\boldsymbol{\mu}_{ML}$, \mathbf{W}_{ML} and $\boldsymbol{\Phi}_{ML}$ represent the maximum likelihood solution corresponding to the original untransformed data, show that $\mathbf{A}\boldsymbol{\mu}_{ML}$, \mathbf{AW}_{ML} and $\mathbf{A}\boldsymbol{\Phi}_{ML}\mathbf{A}^T$ will represent the corresponding maximum likelihood solution for the transformed data set. Finally, show that the form of the model is preserved in two cases:

- (a) \mathbf{A} is a diagonal matrix and $\boldsymbol{\Phi}$ is a diagonal matrix. This corresponds to the case of factor analysis. The transformed $\boldsymbol{\Phi}$ remains diagonal, and hence factor analysis is covariant under component-wise re-scaling of the data variables
 - (b) \mathbf{A} is orthogonal and $\boldsymbol{\Phi}$ is proportional to the unit matrix so that $\boldsymbol{\Phi} = \sigma^2 \mathbf{I}$. This corresponds to probabilistic PCA. The transformed $\boldsymbol{\Phi}$ matrix remains proportional to the unit matrix, and hence probabilistic PCA is covariant under a rotation of the axes of data space, as is the case for conventional PCA
- B. For the data set given [at this link](#), Compute the first 8 kernel principal component basis functions. Use RBF kernel with $\sigma^2 = 0.1$

A. (a)

A is a diagonal matrix

$$A_{D \times D} = \begin{bmatrix} A_{11} & 0 & \cdots & 0 \\ 0 & A_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & A_{DD} \end{bmatrix}$$

$$\boldsymbol{\Phi}_{D \times D} = \begin{bmatrix} \boldsymbol{\Phi}_{11} & 0 & \cdots & 0 \\ 0 & \boldsymbol{\Phi}_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \boldsymbol{\Phi}_{DD} \end{bmatrix}$$

The transformed $\bar{\Phi} \rightarrow A\bar{\Phi}A^T$

$$A\bar{\Phi}A^T = \begin{bmatrix} A_{11} & 0 & \cdots & 0 \\ 0 & A_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & A_{DD} \\ 0 & 0 & \cdots & A_{DD} \end{bmatrix} \begin{bmatrix} \bar{\Phi}_{11} & 0 & \cdots & 0 \\ 0 & \bar{\Phi}_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \bar{\Phi}_{DD} \\ 0 & 0 & \cdots & \bar{\Phi}_{DD} \end{bmatrix} \begin{bmatrix} A_{11} & 0 & \cdots & 0 \\ 0 & A_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & A_{DD} \\ 0 & 0 & \cdots & A_{DD} \end{bmatrix}$$

$$= \begin{bmatrix} A_{11}\bar{\Phi}_{11}A_{11} & 0 & \cdots & 0 \\ 0 & A_{22}\bar{\Phi}_{22}A_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & A_{DD}\bar{\Phi}_{DD}A_{DD} \end{bmatrix}$$

$A\bar{\Phi}A^T$ is diagonal

Therefore, FA is covariant under component-wise rescaling of the data variables.

b.

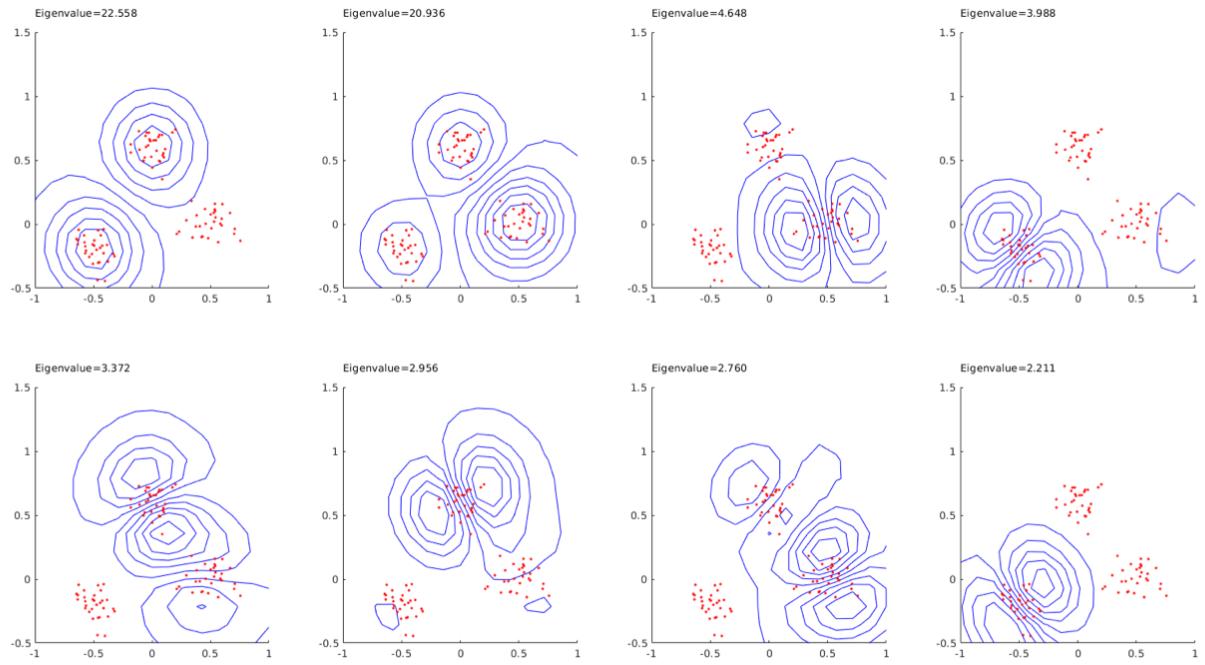
A is orthogonal $AA^T = I$

$$\Phi = \alpha^2 I$$

$$\begin{aligned} A\Phi A^T &= A\alpha^2 I A^T = \alpha^2 A I A^T \\ &= \alpha^2 A A^T \\ &= \alpha^2 I \end{aligned}$$

$A\Phi A^T$ is proportional to the unit matrix.

B.



3 Independent Component Analysis

12.6

- A. Suppose that two variables z_1 and z_2 are independent. Show that the covariance matrix between these variables is diagonal. Now consider two variables y_1 and y_2 in which $-1 \leq y_1 \leq 1$ and $y_2 = y_1^2$. Write down the conditional distribution $p(y_2|y_1)$ and show that this is dependent on y_1 . Now show that the covariance matrix between these two variables is again diagonal. This counter-example shows that zero correlation is not a sufficient condition for independence.
- B. For some noisy observations of a 4D signal, use ICA to reconstruct the signal. The data set can be found [at this link](#).

Figure 12.20

A. (1) z_1 and z_2 are independent.

$$P(z_1, z_2) = P(z_1) P(z_2)$$

$$E(z_1 z_2) = E(z_1) E(z_2)$$

$$\text{corr}(z_1, z_2) = \frac{E[(z_1 - E[z_1])(z_2 - E[z_2])]}{\sqrt{\text{Var}[z_1]} \sqrt{\text{Var}[z_2]}} = 0$$

the covariance matrix between these variables is diagonal.

(2) $-1 \leq y_1 \leq 1$, assume uniform distribution. $P(y_1) = \frac{1}{2}$

$$y_2 = y_1^2, \quad P(y_2, y_1) = \frac{1}{y_1^2}, \quad E[y_1] = 0$$

$$P(y_2|y_1) = \frac{P(y_2, y_1)}{P(y_1)} = \frac{2}{y_1^2}, \text{ dependent on } y_1$$

$$\text{corr}(y_1, y_2) = \frac{E[(y_1 - E[y_1])(y_2 - E[y_2])]}{\sqrt{\text{Var}[y_1]} \sqrt{\text{Var}[y_2]}}$$

$$\text{corr}(y_1, y_2) = \frac{\overbrace{E[y_1(y_2 - E[y_2])]}^{\sqrt{\text{Var}[y_1]}}}{\sqrt{\text{Var}[y_1]}\sqrt{\text{Var}[y_2]}}$$

$$= \frac{\overbrace{E[y_1y_2 - y_1E[y_2]]}^{\sqrt{\text{Var}[y_1]}\sqrt{\text{Var}[y_2]}}}{\sqrt{\text{Var}[y_1]}\sqrt{\text{Var}[y_2]}}$$

y_2 is dependent on y_1

$$y_1 y_2 = y_1 y_1^2 = y_1^3$$

$$y_1 E[y_2] = y_1 y_1^3 = y_1^3$$

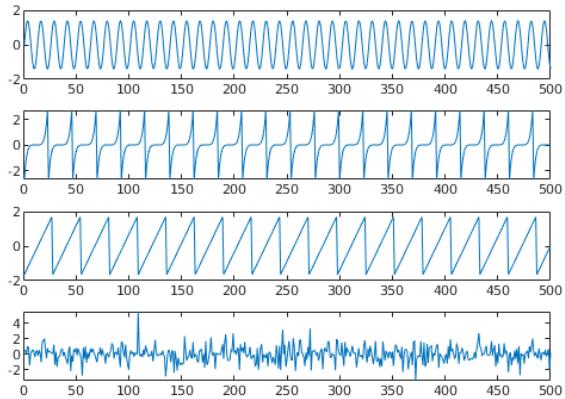
$$E[y_1 y_2 - y_1 E[y_2]] = 0$$

$$\text{corr}(y_1, y_2) = 0$$

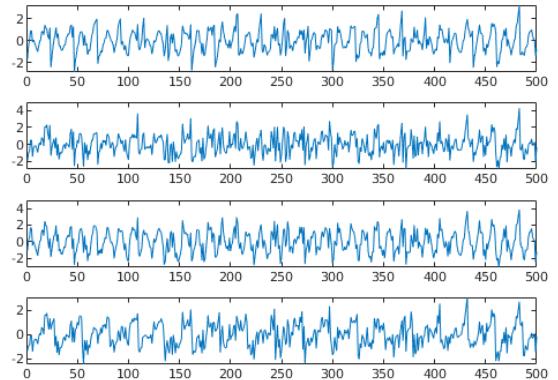
the covariance matrix between y_1, y_2 is again diagonal.

B.

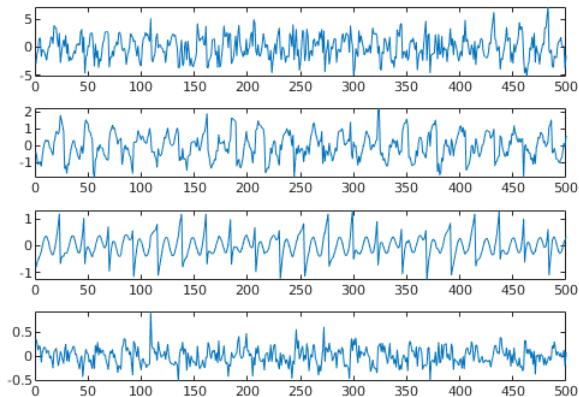
truth



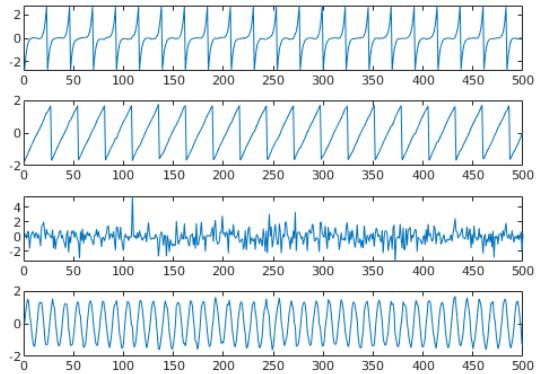
observed signals



PCA estimate



ICA estimate



4 LASSO

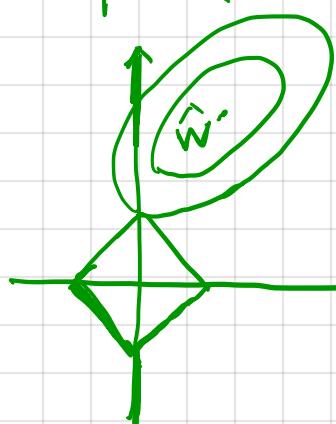
13.5

- A. Recall that the LASSO objective is to minimize $RSS(\beta) + \lambda \sum_j |\beta_j|$ whereas the ridge regression objective is to minimize $RSS(\beta) + \lambda \sum_j \|\beta_j\|_2^2$.
- True or False:** LASSO solutions result in sparsity in the regression coefficients. Explain.
 - True or False:** Ridge Regression solutions result in sparsity in the regression coefficients. Explain
 - True or False:** As we increase λ in LASSO, we expect the number of variables in our solution to increase. Explain.
 - True or False:** It is possible to achieve the least-squares solution using a LASSO objective. Explain
 - True or False:** Given any two LASSO solutions corresponding to λ_1 and λ_2 with $\lambda_2 > \lambda_1$ and the same support for these two solutions, it is possible to write out a closed-form expression for all solutions corresponding to λ with $\lambda_1 < \lambda < \lambda_2$. Explain.
- B. Generate sparse signal \mathbf{w}^* of size $D = 4096$, consisting of 160 randomly placed ± 1 spikes. Next generate a random design matrix \mathbf{X} of size $N \times D$ where $N = 1024$. Finally, generate a noisy observation $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$ where $\epsilon_i = \mathcal{N}(0, 0.01^2)$. Estimate \mathbf{w} from \mathbf{y} and \mathbf{X} . Estimate \mathbf{w} by using LASSO using $\lambda = 0.1\lambda_{max}$. Plot the results and discuss.

p58, slides

A. (a) True

Explain:



$$\text{LASSO: } RSS(\beta) + \lambda \sum_j |\beta_j|$$

We can re-write this as a constrained but smooth objective (a quadratic function with linear constraints)

$$\min_{\beta} RSS(\beta) \quad \text{s.t. } \sum_j |\beta_j| \leq B$$

Contour of the RSS objective function of LASSO

where B is the upper bound on the ℓ_1 -norm of weight

From the theory of constrained optimization, we know that the optimal solution occurs at the point where the lowest level set of the objective function intersects the constraint surface.

It should be geometrically clear that as we relax the constrain B, we "grow" the ℓ_1 ball until it meets the objective.

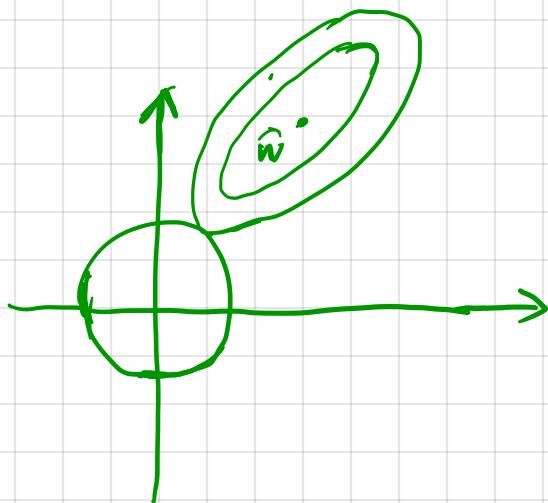
The corners of the ball are more likely to intersect the ellipse than one of the sides, especially in high dimensions, because the corners stick out more.

The corners correspond to sparse solutions, which lie on the coordinate axes.

Therefore, LASSO solutions result in sparsity in the regression coefficients.

(b)

False



$$\min_{\beta} \text{RSS}(\beta) \quad \text{s.t. } \exists j \quad \|\beta_j\|_2^2 \leq B$$

When we grow the ℓ_2 ball, it can intersect the objective at any point, there are no corners, so there is no preference for sparsity.

(C) False

As we increase λ , the solution vector $\hat{w}(\lambda)$ will tend to get sparser.

(d) True

Set $\lambda = 0$

$$\text{LASSO} = \text{RSS}(\beta)$$

(e) True.

Explain:

from the result of regularization path.

$\hat{w}_j(\lambda)$ vs λ

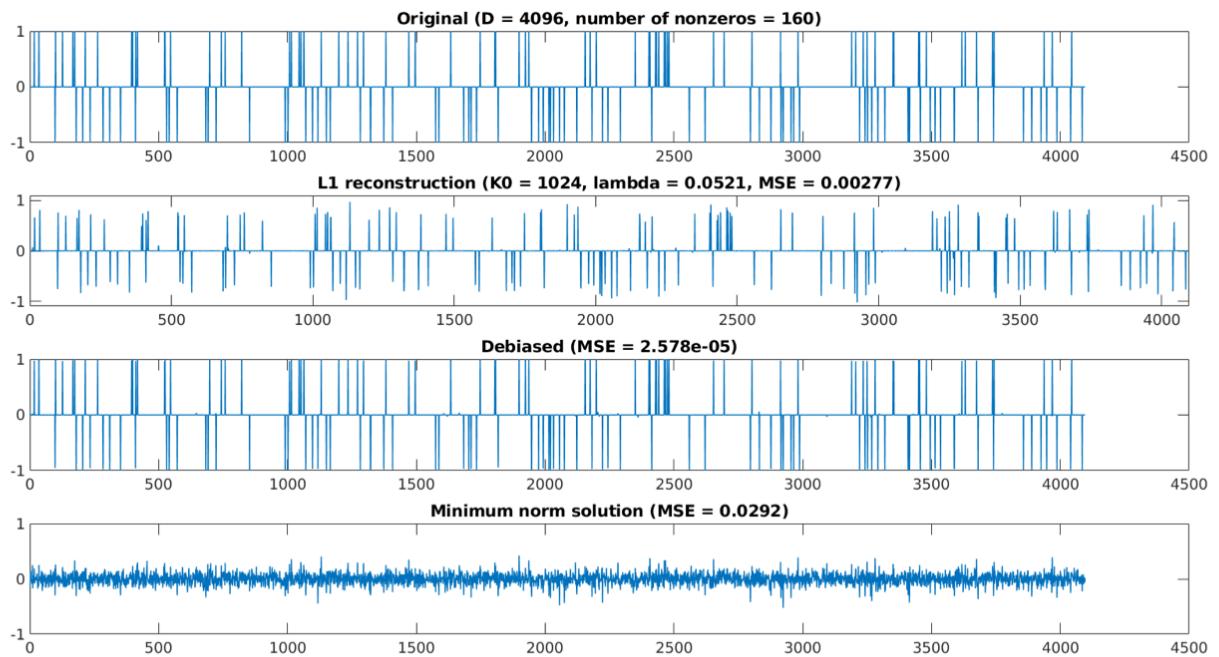


almost flat when $\hat{w}_j(\lambda)$ is small

λ_1, λ_2 has the same support.

it is possible to find a λ^*
 $\lambda_1 < \lambda^* < \lambda_2$, and write a
close-out expression.

B.



5 Automatic relevance determination and compressed sensing

13.7

A. Automatic relevance determination is another method that results in sparse solution.

- (a) Explain (with equations and diagram) how ARD works and why it results in sparse solution.
- (b) How is ARD connected to the MAP estimate?

B. Suppose we have an image which is corrupted in some way, e.g., by having text or scratches sparsely superimposed on top of it. We might want to estimate the underlying clean image. This is called image inpainting. One can use similar techniques for image denoising. One way to address this problem is to use compressed sensing

Write a code for image denoising using compressed sensing. Use it for the image provided [at this link](#).

A. (a)

$$\alpha_j = \frac{1}{\sigma_j^2} \quad \text{weight precision}$$

$$\beta = \frac{1}{\sigma^2} \quad \text{measurement precision}$$

$$P(y|X, w, \beta) = N(y|w^T X, 1/\beta)$$

$$P(w) = N(w|0, A^{-1}), A = \text{diag}(\alpha)$$

The marginal likelihood

$$\begin{aligned} P(y|X, \alpha, \beta) &= \int N(y|Xw, \beta I_N) N(w|0, A) dw \\ &= N(y|0, \beta I_N + XA^{-1}X^T) \\ &= (2\pi)^{-\frac{N}{2}} |C_\alpha|^{-\frac{1}{2}} \exp(-\frac{1}{2} y^T C_\alpha^{-1} y) \end{aligned}$$

$$\text{where } C_\alpha \triangleq \beta^{-1} I_N + X A^T X^T$$

$$l(\alpha, \beta) \triangleq -\frac{1}{2} \log P(y | X, \alpha, \beta)$$

$$\stackrel{\Delta}{=} \log |C_\alpha| + y^T C_\alpha^{-1} y$$

$$\stackrel{\Delta}{=} \log |C_\alpha| + y^T C_\alpha^{-1} y + \sum_j (\alpha + \log \alpha_j - b \alpha_j) + \log \beta - d \beta.$$

Optimize $\hat{\alpha}, \hat{\beta}$.

Once we have estimated $\hat{\alpha}$ and $\hat{\beta}$, we can compute the posterior over the parameter

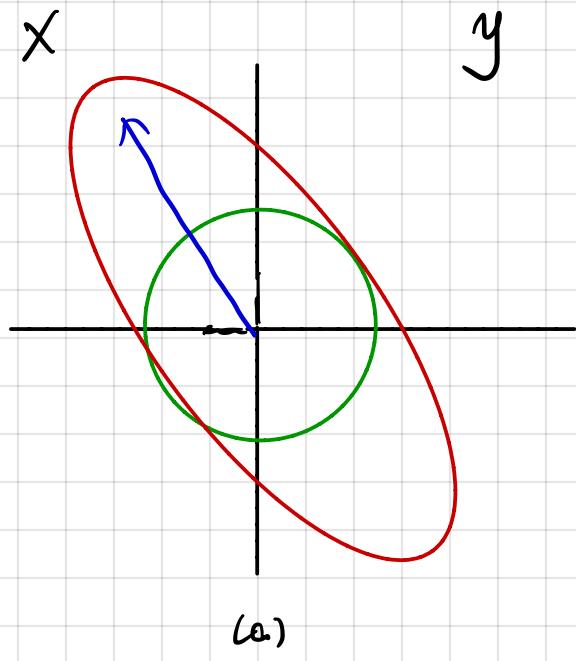
$$P(w | D, \hat{\alpha}, \hat{\beta}) = \mathcal{N}(\mu, \Sigma)$$

$$\Sigma^{-1} = \hat{\beta} X^T X + A$$

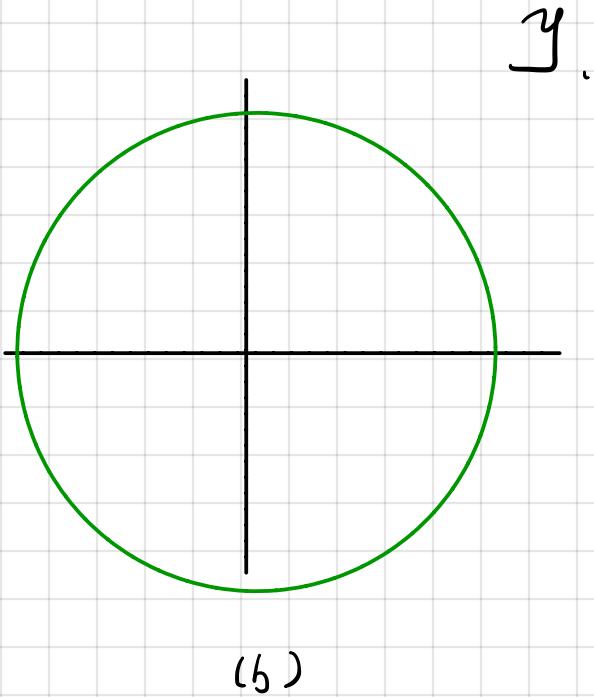
$$\mu = \hat{\beta} \Sigma X^T y$$

The fact that we computer $P(w | D, \hat{\alpha}, \hat{\beta})$, while simultaneously encouraging sparsity.

from Figure 13.20. (Murphy's book)



(a)



(b)

Illustration of why ARD results in sparsity. The vector of input x does not point towards the vector of outputs y , so the feature should be removed. (a) For finite λ , the probability density is spread in directions away from y . (b) when $\lambda = \infty$, the probability density at y is maximized.

The marginal likelihood is

$$P(y|x, \alpha, \beta) = \mathcal{N}(y|\alpha, C)$$

$$C = \frac{1}{\beta} I + \frac{1}{2} X X^T$$

If α is finite, the posterior will be elongated along the direction of X .

However, if $\alpha = \infty$, $C = \frac{1}{\beta} I$, C is spherical.

If $|C|$ is held constant, the latter assigns higher probability density to the observed response vector y , so this is the preferred solution.

(b)

the effective prior $p(w|\hat{\lambda})$ is non-factorial, and further it depends on data D and α^2 .

$$\hat{w}^{\text{ARD}} = \arg \min_w \beta \|y - Xw\|_2^2 + g_{\text{ARD}}(w)$$

$$g_{\text{ARD}}(w) \triangleq \min_{\alpha > 0} \sum_j \alpha_j w_j^2 + \log |\mathcal{C}_2|$$

ARD can be viewed as the MAP estimation problem.

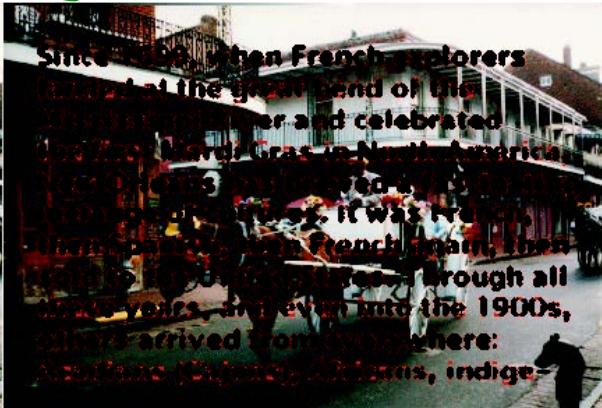
MAP estimation with non-factorial priors is strictly better than MAP estimation with any possible factorial prior.

B.

original



remove red text



then use OMP to recover the image.



remove the red text

RGB, if $R \geq 180$, $G \leq 60$, $B \leq 60$,

the $R=0$, $G=0$, $B=0$

the I use OMP to recover the image.

the code is from: [github chongyangtao/color-Image-Inpainting](https://github.com/chongyangtao/color-Image-Inpainting).